

How well can we estimate a sparse vector?

Emmanuel J. Candès and Mark A. Davenport*

March 4, 2013

Abstract

The estimation of a sparse vector in the linear model is a fundamental problem in signal processing, statistics, and compressive sensing. This paper establishes a lower bound on the mean-squared error, which holds regardless of the sensing/design matrix being used and regardless of the estimation procedure. This lower bound very nearly matches the known upper bound one gets by taking a random projection of the sparse vector followed by an ℓ_1 estimation procedure such as the Dantzig selector. In this sense, compressive sensing techniques cannot essentially be improved.

Keywords: Compressive sensing, sparse estimation, sparse linear regression, minimax lower bounds, Fano's inequality, matrix Bernstein inequality

1 Introduction

The estimation of a sparse vector from noisy observations is a fundamental problem in signal processing and statistics, and lies at the heart of the growing field of compressive sensing [4, 5, 8]. At its most basic level, we are interested in accurately estimating a vector $\mathbf{x} \in \mathbb{R}^n$ that has at most k non-zeros from a set of noisy linear measurements

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{z}, \quad (1)$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$. We are often interested in the underdetermined setting where m may be much smaller than n . In general, one would not expect to be able to accurately recover \mathbf{x} when $m < n$ since there are more unknowns than observations. However it is by now well-known that by exploiting sparsity, it is possible to accurately estimate \mathbf{x} .

As an example, consider what is known concerning ℓ_1 minimization techniques, which are among the most powerful and well-understood with respect to their performance in noise. Specifically, if we suppose that the entries of the matrix \mathbf{A} are i.i.d. $\mathcal{N}(0, 1/n)$, then one can show that for any $\mathbf{x} \in \Sigma_k := \{\mathbf{x} : \|\mathbf{x}\|_0 \leq k\}$, ℓ_1 minimization techniques such as the Lasso or the Dantzig selector produce a recovery $\hat{\mathbf{x}}$ such that

$$\frac{1}{n} \|\hat{\mathbf{x}} - \mathbf{x}\|_2^2 \leq C_0 \frac{k\sigma^2}{m} \log n \quad (2)$$

holds with high probability provided that $m = \Omega(k \log(n/k))$ [6]. We refer to [3] and [9] for further results.

*E.J.C. is with the Departments of Mathematics and Statistics, Stanford University, Stanford, CA 94035. M.A.D. is with the School of Electrical and Computer Engineering, Georgia Institute of Technology, Atlanta, GA 30332. This work has been partially supported by NSF grant DMS-1004718, the Waterman Award from NSF, ONR grant N00014-10-1-0599 and a grant from AFOSR. M.A.D. is the corresponding author. Email: mdav@gatech.edu. Phone: (404)894-2881. Fax: (404)894-8363.

1.1 Criticism

A noteworthy aspect of the bound in (2) is that the recovery error increases linearly as we decrease m , and thus we pay a penalty for taking a small number of measurements. Although this effect is sometimes cited as a drawback of the compressive sensing framework, it should not be surprising — we fully expect that if each measurement has a constant SNR, then taking more measurements should reduce our estimation error.

However, there is another somewhat more troubling aspect of (2). Specifically, by filling the rows of \mathbf{A} with i.i.d. random variables, we are ensuring that our “sensing vectors” are almost orthogonal to our signal of interest, leading to a tremendous SNR loss. To quantify this loss, suppose that we had access to an oracle that knows *a priori* the locations of the nonzero entries of \mathbf{x} and could instead construct \mathbf{A} with vectors localized to the support of \mathbf{x} . For example, if m is an integer multiple of k then we could simply measure sample each coefficient directly m/k times and then average these samples. One can check that this procedure would yield an estimate obeying

$$\mathbb{E} \left[\frac{1}{n} \|\hat{\mathbf{x}} - \mathbf{x}\|_2^2 \right] = \left(\frac{k\sigma^2}{m} \right) \left(\frac{k}{n} \right). \quad (3)$$

Thus, the performance in (2) is worse than what would be possible with an oracle by a factor of $(n/k) \log n$. When k is small, this is very large! Of course, we won’t have access to an oracle in practice, but the substantial difference between (2) and (3) naturally leads one to question whether (2) can be improved upon.

1.2 Can we do better?

In this paper we will approach this question from the viewpoint of compressive sensing and/or of experimental design. Specifically, we assume that we are free to choose *both* the matrix \mathbf{A} and the sparse recovery algorithm. Our results will have implications for the case where \mathbf{A} is determined by factors beyond our control, but our primary interest will be in considering the performance obtained by the best possible choice of \mathbf{A} . In this setting, our fundamental question is:

Can we ever hope to do better than (2)? Is there a more intelligent choice for the matrix \mathbf{A} ? Is there a more effective recovery algorithm?

In this paper we show that the answer is *no*, and that there exists no choice of \mathbf{A} or recovery algorithm that can significantly improve upon the guarantee in (2). Specifically, we consider the worst-case error over all $\mathbf{x} \in \Sigma_k$, i.e.,

$$M^*(\mathbf{A}) = \inf_{\hat{\mathbf{x}}} \sup_{\mathbf{x} \in \Sigma_k} \mathbb{E} \left[\frac{1}{n} \|\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x}\|_2^2 \right]. \quad (4)$$

Our main result consists of the following bound, which establishes a fundamental limit on the minimax risk which holds for any matrix \mathbf{A} and any possible recovery algorithm.

Theorem 1. *Suppose that we observe $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{z}$ where \mathbf{x} is a k -sparse vector, \mathbf{A} is an $m \times n$ matrix with $m \geq k$, and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$. Then there exists a constant $C_1 > 0$ such that for all \mathbf{A} ,*

$$M^*(\mathbf{A}) \geq C_1 \frac{k\sigma^2}{\|\mathbf{A}\|_F^2} \log(n/k). \quad (5)$$

We also have that for all \mathbf{A}

$$M^*(\mathbf{A}) \geq \frac{k\sigma^2}{\|\mathbf{A}\|_F^2}. \quad (6)$$

This theorem says that *there is no \mathbf{A} and no recovery algorithm* that does fundamentally better than the Dantzig selector (2) up to a constant¹; that is, ignoring the difference in the factors $\log n/k$ and $\log n$. In this sense, the results of compressive sensing are at the limit.

Although the noise model in (1) is fairly common, in some settings (such as the estimation of a signal transmitted over a noisy channel) it is more natural to consider noise that has been added directly to the signal prior to the acquisition of the measurements. In this case we can directly apply Theorem 1 to obtain the following corollary.

Corollary 1. *Suppose that we observe $\mathbf{y} = \mathbf{A}(\mathbf{x} + \mathbf{w})$ where \mathbf{x} is a k -sparse vector, \mathbf{A} is an $m \times n$ matrix with $k \leq m \leq n$, and $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$. Then for all \mathbf{A}*

$$M^*(\mathbf{A}) \geq C_1 \frac{k\sigma^2}{m} \log(n/k) \quad \text{and} \quad M^*(\mathbf{A}) \geq \frac{k\sigma^2}{m}. \quad (7)$$

Proof. We assume that \mathbf{A} has rank $m' \leq m$. Let $\mathbf{U}\Sigma\mathbf{V}^*$ be the reduced SVD of \mathbf{A} , where \mathbf{U} is $m \times m'$, Σ is $m' \times m'$, and \mathbf{V} is $n \times m'$. Applying the matrix $\Sigma^{-1}\mathbf{U}^*$ to \mathbf{y} preserves all the information about \mathbf{x} , and so we can equivalently assume that the data is given by

$$\mathbf{y}' = \Sigma^{-1}\mathbf{U}^*\mathbf{y} = \mathbf{V}^*\mathbf{x} + \mathbf{V}^*\mathbf{w}. \quad (8)$$

Note that $\mathbf{V}^*\mathbf{w}$ is a Gaussian vector with covariance matrix $\sigma^2 \mathbf{V}^*\mathbf{V} = \sigma^2 \mathbf{I}$. Moreover, \mathbf{V}^* has unit-norm rows, so that $\|\mathbf{V}^*\|_F \leq m' \leq m$. We then apply Theorem 1 to establish (7). \square

The intuition behind this result is that when noise is added to the measurements, we can boost the SNR by rescaling \mathbf{A} to have higher norm. When we instead add noise to the signal, the noise is also scaled by \mathbf{A} , and so no matter how \mathbf{A} is designed there will always be a penalty of $1/m$.

1.3 Related work

There have been a number of prior works that have established lower bounds on $M^*(\mathbf{A})$ or related quantities under varying assumptions [1, 13–17, 19]. In [1, 17], techniques from information theory similar to the ones that we use below are used to establish rather general lower bounds under the assumption that the entries of \mathbf{x} are generated i.i.d. according to some distribution. For an appropriate choice of distribution, \mathbf{x} will be approximately sparse and [1, 17] will yield asymptotic lower bounds of a similar flavor to ours.

The prior work most closely related to our results is that of Ye and Zhang [19] and Raskutti, Wainwright, and Yu [15]. In [19] Ye and Zhang establish a bound similar to (5) in Theorem 1. While the resulting bounds are substantially the same, the bounds in [19] hold only in the asymptotic regime where $k \rightarrow \infty$, $n \rightarrow \infty$, and $\frac{k}{n} \rightarrow 0$, whereas our bounds hold for arbitrary finite values of k and n , including the case where k is relatively large compared to n . In [15] Raskutti et al. reach a somewhat similar conclusion to our Theorem 1 via a similar argument, but where it is assumed that \mathbf{A} satisfies $\|\mathbf{A}\mathbf{x}\|_2^2 \leq (1 + \delta) \|\mathbf{x}\|_2^2$ for all $\mathbf{x} \in \Sigma_{2k}$, (i.e., the upper bound of the *restricted isometry*

¹Our analysis shows that asymptotically C_1 can be taken as 1/128. We have made no effort to optimize this constant, and it is probably far from sharp. This is why we give the simpler bound (6) which is proven by considering the error we would incur even if we knew the support of \mathbf{x} *a priori*. However, our main result is (5). We leave the calculation of an improved constant to future work.

property (RIP)). In this case the authors show that²

$$M^*(\mathbf{A}) \geq C \frac{k\sigma^2}{(1+\delta)n} \log(n/k).$$

Our primary aim, however, is to challenge the use of the RIP and/or random matrices and to determine whether we can do better via a different choice in \mathbf{A} . Our approach relies on standard tools from information theory such as Fano's inequality, and as such is very similar in spirit to the approaches in [1, 15, 17]. The proof of Theorem 1 begins by following a similar path to that taken in [15]. As in the results of [15], we rely on the construction of a packing set of sparse vectors. However, we place no assumptions whatsoever on the matrix \mathbf{A} . To do this we must instead consider a random construction of this set, allowing us to apply the recently established matrix-version of Bernstein's inequality due to Ahlswede and Winter [2] to bound the empirical covariance matrix of the packing set. Our analysis is divided into two parts. In Section 2 we provide the proof of Theorem 1, and in Section 3 we provide the construction of the necessary packing set.

1.4 Notation

We now provide a brief summary of the notations used throughout the paper. If \mathbf{A} is an $m \times n$ matrix and $T \subset \{1, \dots, n\}$, then \mathbf{A}_T denotes the $m \times |T|$ submatrix with columns indexed by T . Similarly, for a vector $\mathbf{x} \in \mathbb{R}^n$ we let $\mathbf{x}|_T$ denote the restriction of \mathbf{x} to T . We will use $\|\mathbf{x}\|_p$ to denote the standard ℓ_p norm of a vector, and for a matrix \mathbf{A} , we will use $\|\mathbf{A}\|$ and $\|\mathbf{A}\|_F$ to denote the operator and Frobenius norms respectively.

2 Proof of Main Result

In this section we establish the lower bound (5) in Theorem 1. The proof of (6) is provided in the Appendix. In the proofs of both (5) and (6), we will assume that $\sigma = 1$ since the proof for arbitrary σ follows by a simple rescaling. To obtain the bound in (5) we begin by following a similar course as in [15]. Specifically, we will suppose that \mathbf{x} is distributed uniformly on a finite set of points $\mathcal{X} \subset \Sigma_k$, where \mathcal{X} is constructed so that the elements of \mathcal{X} are well separated. This allows us to apply the following lemma which follows from Fano's inequality combined with the convexity of the Kullback-Leibler (KL) divergence. We provide a proof of the lemma in the Appendix.

Lemma 1. *Consider the measurement model where $\mathbf{y} = \mathbf{Ax} + \mathbf{z}$ with $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Suppose that there exists set of points $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^{|\mathcal{X}|} \subset \Sigma_k$ such that for any $\mathbf{x}_i, \mathbf{x}_j \in \mathcal{X}$, $\|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \geq 8nM^*(\mathbf{A})$, where $M^*(\mathbf{A})$ is defined as in (4). Then*

$$\frac{1}{2} \log |\mathcal{X}| - 1 \leq \frac{1}{2|\mathcal{X}|^2} \sum_{i,j=1}^{|\mathcal{X}|} \|\mathbf{Ax}_i - \mathbf{Ax}_j\|_2^2. \quad (9)$$

²Note that it is possible to remove the assumption that \mathbf{A} satisfies the upper bound of the RIP, but with a rather unsatisfying result. Specifically, for an arbitrary matrix \mathbf{A} with a fixed Frobenius norm, we have that $\|\mathbf{A}\|_2^2 \leq \|\mathbf{A}\|_F^2$, so that $(1+\delta) \leq \|\mathbf{A}\|_F^2$. This bound can be shown to be tight by considering a matrix \mathbf{A} with only one nonzero column. However, applying this bound underestimates $M^*(\mathbf{A})$ by a factor of n . Of course, the bounds coincide for “good” matrices (such as random matrices) which will have a significantly smaller value of δ [14]. However, the random matrix framework is precisely that which we wish to challenge.

By taking the set \mathcal{X} in Lemma 2 below and rescaling these points by $4\sqrt{nM^*(\mathbf{A})}$, we have that there exists a set \mathcal{X} satisfying the assumptions of Lemma 1 with

$$|\mathcal{X}| = (n/k)^{k/4},$$

and hence from (9) we obtain

$$\frac{k}{4} \log(n/k) - 2 \leq \frac{1}{|\mathcal{X}|^2} \sum_{i,j=1}^{|\mathcal{X}|} \|\mathbf{A}\mathbf{x}_i - \mathbf{A}\mathbf{x}_j\|_2^2 = \text{Tr} \left(\mathbf{A}^* \mathbf{A} \left(\frac{1}{|\mathcal{X}|^2} \sum_{i,j=1}^{|\mathcal{X}|} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^* \right) \right). \quad (10)$$

If we set

$$\boldsymbol{\mu} = \frac{1}{|\mathcal{X}|} \sum_{i=1}^{|\mathcal{X}|} \mathbf{x}_i \quad \text{and} \quad \mathbf{Q} = \frac{1}{|\mathcal{X}|} \sum_{i=1}^{|\mathcal{X}|} \mathbf{x}_i \mathbf{x}_i^*,$$

then one can show that

$$\frac{1}{|\mathcal{X}|^2} \sum_{i,j=1}^{|\mathcal{X}|} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^* = 2(\mathbf{Q} - \boldsymbol{\mu}\boldsymbol{\mu}^*).$$

Thus, we can bound (10) by

$$2 \text{Tr}(\mathbf{A}^* \mathbf{A} (\mathbf{Q} - \boldsymbol{\mu}\boldsymbol{\mu}^*)) \leq 2 \text{Tr}(\mathbf{A}^* \mathbf{A} \mathbf{Q}),$$

where the inequality follows since $\text{Tr}(\mathbf{A}^* \mathbf{A} \boldsymbol{\mu}\boldsymbol{\mu}^*) = \|\mathbf{A}\boldsymbol{\mu}\|_2^2 \geq 0$. Moreover, since $\mathbf{A}^* \mathbf{A}$ and \mathbf{Q} are positive semidefinite,

$$\text{Tr}(\mathbf{A}^* \mathbf{A} \mathbf{Q}) \leq \text{Tr}(\mathbf{A}^* \mathbf{A}) \|\mathbf{Q}\| = \|\mathbf{A}\|_F^2 \|\mathbf{Q}\|.$$

Combining this with (10) and applying Lemma 2 to bound the norm of \mathbf{Q} — recalling that it has been appropriately rescaled — we obtain

$$\frac{k}{4} \log(n/k) - 2 \leq (1 + \beta) 32 M^*(\mathbf{A}) \|\mathbf{A}\|_F^2,$$

where β is a constant that can be arbitrarily close to 0. This yields the desired result.

3 Packing Set Construction

We now return to the problem of constructing the packing set \mathcal{X} . As noted above, our construction exploits the following matrix Bernstein inequality of Ahlswede and Winter [2]. See also [18].

Theorem 2 (Matrix Bernstein Inequality). *Let $\{\mathbf{X}_i\}$ be a finite sequence of independent zero-mean random self-adjoint matrices of dimension $n \times n$. Suppose that $\|\mathbf{X}_i\| \leq 1$ almost surely for all i and set $\rho^2 = \sum_i \|\mathbb{E}[\mathbf{X}_i^2]\|$. Then for all $t \in [0, 2\rho^2]$,*

$$\mathbb{P} \left[\left\| \sum_i \mathbf{X}_i \right\| \geq t \right] \leq 2n \exp \left(-\frac{t^2}{4\rho^2} \right). \quad (11)$$

We construct the set \mathcal{X} by choosing points at random, which allows us to apply Theorem 2 to establish a bound on the empirical covariance matrix. In bounding the size of \mathcal{X} we follow a similar course as in [15] and rely on techniques from [11].

Lemma 2. Let n and k be given, and suppose for simplicity that k is even and $k < n/2$. There exists a set $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^{|\mathcal{X}|} \subset \Sigma_k$ of size

$$|\mathcal{X}| = (n/k)^{k/4} \quad (12)$$

such that

(i) $\|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \geq 1/2$ for all $\mathbf{x}_i, \mathbf{x}_j \in \mathcal{X}$ with $i \neq j$; and

(ii) $\left\| \frac{1}{|\mathcal{X}|} \sum_{i=1}^{|\mathcal{X}|} \mathbf{x}_i \mathbf{x}_i^* - \frac{1}{n} \mathbf{I} \right\| \leq \beta/n$,

where β can be made arbitrarily close to 0 as $n \rightarrow \infty$.

Proof. We will show that such a set \mathcal{X} exists via the probabilistic method. Specifically, we will show that if we draw $|\mathcal{X}|$ independent k -sparse vectors at random, then the set will satisfy both (i) and (ii) with probability strictly greater than 0. We will begin by considering the set

$$\mathcal{U} = \left\{ \mathbf{x} \in \left\{ 0, +\sqrt{1/k}, -\sqrt{1/k} \right\}^n : \|\mathbf{x}\|_0 = k \right\}.$$

Clearly, $|\mathcal{U}| = \binom{n}{k} 2^k$. Next, note that for all $\mathbf{x}, \mathbf{x}' \in \mathcal{U}$, $\frac{1}{k} \|\mathbf{x}' - \mathbf{x}\|_0 \leq \|\mathbf{x}' - \mathbf{x}\|_2^2$, and thus if $\|\mathbf{x}' - \mathbf{x}\|_2^2 \leq 1/2$ then $\|\mathbf{x}' - \mathbf{x}\|_0 \leq k/2$. From this we observe that for any fixed $\mathbf{x} \in \mathcal{U}$,

$$\left| \left\{ \mathbf{x}' \in \mathcal{U} : \|\mathbf{x}' - \mathbf{x}\|_2^2 \leq 1/2 \right\} \right| \leq \left| \left\{ \mathbf{x}' \in \mathcal{U} : \|\mathbf{x}' - \mathbf{x}\|_0 \leq k/2 \right\} \right| \leq \binom{n}{k/2} 3^{k/2}.$$

Suppose that we construct \mathcal{X} by picking elements of \mathcal{U} uniformly at random. When adding the j^{th} point to \mathcal{X} , the probability that \mathbf{x}_j violates (i) with respect to the previously added points is bounded by

$$\frac{(j-1) \binom{n}{k/2} 3^{k/2}}{\binom{n}{k} 2^k}.$$

Thus, using the union bound, we can bound the total probability that \mathcal{X} will fail to satisfy (i), denoted P_1 , by

$$P_1 \leq \sum_{j=1}^{|\mathcal{X}|} \frac{(j-1) \binom{n}{k/2} 3^{k/2}}{\binom{n}{k} 2^k} \leq \frac{|\mathcal{X}|^2 \binom{n}{k/2}}{2 \binom{n}{k}} \left(\frac{\sqrt{3}}{2} \right)^k.$$

Next, observe that

$$\frac{\binom{n}{k}}{\binom{n}{k/2}} = \frac{(k/2)!(n-k/2)!}{k!(n-k)!} = \prod_{i=1}^{k/2} \frac{n-k+i}{k/2+i} \geq \left(\frac{n-k+k/2}{k/2+k/2} \right)^{k/2} = \left(\frac{n}{k} - \frac{1}{2} \right)^{k/2},$$

where the inequality follows since $(n-k+i)/(k/2+i)$ is decreasing as a function of i provided that $n-k > k/2$. Also,

$$\left(\frac{n}{k} \right)^{k/2} \left(\frac{\sqrt{3}}{2} \right)^k = \left(\frac{3n}{4k} \right)^{k/2} \leq \left(\frac{n}{k} - \frac{1}{2} \right)^{k/2}$$

with the proviso $k \leq n/2$. Thus, for $|\mathcal{X}|$ of size given in (12),

$$P_1 \leq \frac{1}{2} \left(\frac{n}{k} \right)^{k/2} \frac{\binom{n}{k/2}}{\binom{n}{k}} \left(\frac{\sqrt{3}}{2} \right)^k \leq \frac{1}{2} \left(\frac{n}{k} - \frac{1}{2} \right)^{k/2} \frac{\binom{n}{k/2}}{\binom{n}{k}} \leq \frac{1}{2} \frac{\binom{n}{k}}{\binom{n}{k/2}} \frac{\binom{n}{k/2}}{\binom{n}{k}} \leq \frac{1}{2}. \quad (13)$$

Next, we consider (ii). We begin by letting

$$\mathbf{X}_i = \mathbf{x}_i \mathbf{x}_i^* - \frac{\mathbf{I}}{n}.$$

Since \mathbf{x}_i is drawn uniformly at random from \mathcal{U} , it is straightforward to show that $\|\mathbf{X}_i\| \leq 1$ and that $\mathbb{E}[\mathbf{x}_i \mathbf{x}_i^*] = \mathbf{I}/n$, which implies that $\mathbb{E}[\mathbf{X}_i] = 0$. Moreover,

$$\mathbb{E}[\mathbf{X}_i^2] = \mathbb{E}[(\mathbf{x}_i \mathbf{x}_i^*)^2] - \left(\frac{1}{n} \mathbf{I}\right)^2 = \frac{(n-1)}{n^2} \mathbf{I}.$$

Thus we obtain $\rho^2 = \sum_{i=1}^{|\mathcal{X}|} \|\mathbb{E}[\mathbf{X}_i^2]\| = |\mathcal{X}|(n-1)/n^2 \leq |\mathcal{X}|/n$. Hence, we can apply Theorem 2 to obtain

$$\mathbb{P}\left[\left\|\sum_{i=1}^{|\mathcal{X}|} \mathbf{X}_i\right\| \geq t\right] \leq 2n \exp\left(-\frac{t^2 n}{4|\mathcal{X}|}\right).$$

Setting $t = |\mathcal{X}| \beta/n$, this reduces to show that the probability that \mathcal{X} will fail to satisfy (ii), denoted P_2 , is bounded by

$$P_2 \leq 2n \exp\left(-\frac{\beta^2 |\mathcal{X}|}{4n}\right).$$

For the lemma to hold we require that $P_1 + P_2 < 1$, and since $P_1 < \frac{1}{2}$ it is sufficient to show that $P_2 < \frac{1}{2}$. This will occur provided that

$$\beta^2 > \frac{4n \log(4n)}{|\mathcal{X}|}.$$

Since $|\mathcal{X}| = \Theta((n/k)^k)$, β can be made arbitrarily small as $n \rightarrow \infty$. \square

Appendix

Proof of (6) in Theorem 1. We begin by noting that

$$M^*(\mathbf{A}) = \inf_{\widehat{\mathbf{x}}} \sup_{T:|T| \leq k} \sup_{\mathbf{x}: \text{supp}(\mathbf{x})=T} \mathbb{E}\left[\frac{1}{n} \|\widehat{\mathbf{x}}(\mathbf{y}) - \mathbf{x}\|_2^2\right] \geq \sup_{T:|T| \leq k} \inf_{\widehat{\mathbf{x}}} \sup_{\mathbf{x}: \text{supp}(\mathbf{x})=T} \mathbb{E}\left[\frac{1}{n} \|\widehat{\mathbf{x}}(\mathbf{y}) - \mathbf{x}\|_2^2\right].$$

Thus for the moment we restrict our attention to the subproblem of bounding

$$M^*(\mathbf{A}_T) = \inf_{\widehat{\mathbf{x}}} \sup_{\mathbf{x}: \text{supp}(\mathbf{x})=T} \mathbb{E}\left[\frac{1}{n} \|\widehat{\mathbf{x}}(\mathbf{y}) - \mathbf{x}\|_2^2\right] = \inf_{\widehat{\mathbf{x}}} \sup_{\mathbf{x} \in \mathbb{R}^k} \mathbb{E}\left[\frac{1}{n} \|\widehat{\mathbf{x}}(\mathbf{A}_T \mathbf{x} + \mathbf{z}) - \mathbf{x}\|_2^2\right], \quad (14)$$

where $\widehat{\mathbf{x}}(\cdot)$ takes values in \mathbb{R}^k . The last equality of (14) follows since if $\text{supp}(\mathbf{x}) = T$ then

$$\|\widehat{\mathbf{x}}(\mathbf{y}) - \mathbf{x}\|_2^2 = \|\widehat{\mathbf{x}}(\mathbf{y})|_T - \mathbf{x}|_T\|_2^2 + \|\widehat{\mathbf{x}}(\mathbf{y})|_{T^c}\|_2^2,$$

so that the risk can always be decreased by setting $\widehat{\mathbf{x}}(\mathbf{y})|_{T^c} = 0$. This subproblem (14) has a well-known solution (see Exercise 5.8 on pp. 403 of [12]). Specifically, let $\lambda_i(\mathbf{A}_T^* \mathbf{A}_T)$ denote the eigenvalues of the matrix $\mathbf{A}_T^* \mathbf{A}_T$. Then

$$M^*(\mathbf{A}_T) = \frac{1}{n} \sum_{i=1}^k \frac{1}{\lambda_i(\mathbf{A}_T^* \mathbf{A}_T)}. \quad (15)$$

Thus we obtain

$$M^*(\mathbf{A}) \geq \sup_{T:|T| \leq k} M^*(\mathbf{A}_T) = \sup_{T:|T| \leq k} \frac{1}{n} \sum_{i=1}^k \frac{1}{\lambda_i(\mathbf{A}_T^* \mathbf{A}_T)}. \quad (16)$$

Note that if there exists a subset T for which \mathbf{A}_T is not full rank, then at least one of the eigenvalues $\lambda_i(\mathbf{A}_T^* \mathbf{A}_T)$ will vanish and the minimax risk will be unbounded. This also shows that the minimax risk is always unbounded when $m < k$.

Thus, we now assume that \mathbf{A}_T is full rank for any choice of T . Since $f(x) = 1/x$ is a convex function for $x > 0$, we have that

$$\sum_{i=1}^k \frac{1}{\lambda_i(\mathbf{A}_T^* \mathbf{A}_T)} \geq \frac{k^2}{\sum_{i=1}^k \lambda_i(\mathbf{A}_T^* \mathbf{A}_T)} = \frac{k^2}{\|\mathbf{A}_T\|_F^2}.$$

Since there always exists a set of k columns T_0 such that $\|\mathbf{A}_{T_0}\|_F^2 \leq (k/n) \|\mathbf{A}\|_F^2$, (16) reduces to yield the desired result. \square

Proof of Lemma 1. To begin, note that if \mathbf{x} is uniformly distributed on the set of points in \mathcal{X} , then there exists an estimator $\hat{\mathbf{x}}(\mathbf{y})$ such that

$$\mathbb{E}_{\mathbf{x}, \mathbf{z}} \left[\frac{1}{n} \|\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x}\|_2^2 \right] \leq M^*(\mathbf{A}), \quad (17)$$

where the expectation is now taken with respect to both the signal and the noise. We next consider the problem of deciding which $\mathbf{x}_i \in \mathcal{X}$ generated the observations \mathbf{y} . Towards this end, set

$$T(\hat{\mathbf{x}}(\mathbf{y})) = \arg \min_{\mathbf{x}_i \in \mathcal{X}} \|\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x}_i\|_2.$$

Define $P_e = \mathbb{P}[T(\hat{\mathbf{x}}(\mathbf{y})) \neq \mathbf{x}]$. From Fano's inequality [7] we have that

$$H(\mathbf{x}|\mathbf{y}) \leq 1 + P_e \log |\mathcal{X}|. \quad (18)$$

We now aim to bound P_e . We begin by noting that for any $\mathbf{x}_i \in \mathcal{X}$ and any $\hat{\mathbf{x}}(\mathbf{y})$, $T(\hat{\mathbf{x}}(\mathbf{y})) \neq \mathbf{x}_i$ if and only if there exists an $\mathbf{x}_j \in \mathcal{X}$ with $j \neq i$ such that

$$\|\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x}_i\|_2 \geq \|\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x}_j\|_2 \geq \|\mathbf{x}_i - \mathbf{x}_j\|_2 - \|\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x}_i\|_2.$$

This would imply that

$$2 \|\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x}_i\|_2 \geq \|\mathbf{x}_i - \mathbf{x}_j\|_2 \geq \sqrt{8nM^*(\mathbf{A})}.$$

Thus, we can bound P_e using Markov's inequality as follows:

$$P_e \leq \mathbb{P} \left[\|\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x}_i\|_2^2 \geq 8nM^*(\mathbf{A})/4 \right] \leq \frac{\mathbb{E}_{\mathbf{x}, \mathbf{z}} [\|\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x}_i\|_2^2]}{2nM^*(\mathbf{A})} \leq \frac{nM^*(\mathbf{A})}{2nM^*(\mathbf{A})} = \frac{1}{2}.$$

Combining this with (18) and the fact that $H(\mathbf{x}) = \log |\mathcal{X}|$, we obtain

$$I(\mathbf{x}, \mathbf{y}) = H(\mathbf{x}) - H(\mathbf{x}|\mathbf{y}) \geq \frac{1}{2} \log |\mathcal{X}| - 1.$$

From the convexity of KL divergence (see [10] for details), we have that

$$I(\mathbf{x}, \mathbf{y}) \leq \frac{1}{|\mathcal{X}|^2} \sum_{j,k=1}^{|\mathcal{X}|} D(\mathcal{P}_i, \mathcal{P}_j),$$

where $D(\mathcal{P}_i, \mathcal{P}_j)$ represents the KL divergence from \mathcal{P}_i to \mathcal{P}_j where \mathcal{P}_i denotes the distribution of \mathbf{y} conditioned on $\mathbf{x} = \mathbf{x}_i$. Since $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, \mathcal{P}_i is simply given by $\mathcal{N}(\mathbf{A}\mathbf{x}_i, \mathbf{I})$. Standard calculations demonstrate that $D(\mathcal{P}_i, \mathcal{P}_j) = \frac{1}{2} \|\mathbf{A}\mathbf{x}_i - \mathbf{A}\mathbf{x}_j\|_2^2$, establishing (9). \square

References

- [1] S. Aeron, V. Saligrama, and M. Zhao. Information theoretic bounds for compressed sensing. *IEEE Trans. Inform. Theory*, 56(10):5111–5130, 2010.
- [2] R. Ahlswede and A. Winter. Strong converse for identification via quantum channels. *IEEE Trans. Inform. Theory*, 48(3):569–579, 2002.
- [3] P. Bickel, Y. Ritov, and A. Tsybakov. Simultaneous analysis of Lasso and Dantzig selector. *Ann. Stat.*, 37(4):1705–1732, 2009.
- [4] E. Candès, J. Romberg, and T. Tao. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE Trans. Inform. Theory*, 52(2):489–509, 2006.
- [5] E. Candès and T. Tao. Near optimal signal recovery from random projections: Universal encoding strategies? *IEEE Trans. Inform. Theory*, 52(12):5406–5425, 2006.
- [6] E. Candès and T. Tao. The Dantzig selector: Statistical estimation when p is much larger than n . *Ann. Stat.*, 35(6):2313–2351, 2007.
- [7] T. Cover and J. Thomas. *Elements of Information Theory*. Wiley-Interscience, New York, NY, 1991.
- [8] D. Donoho. Compressed sensing. *IEEE Trans. Inform. Theory*, 52(4):1289–1306, 2006.
- [9] D. Donoho, I. Johnstone, A. Maleki, and A. Montanari. Compressed sensing over ℓ_p balls: Minimax mean square error. *Arxiv preprint arXiv:1103.1943*, 2011.
- [10] T. Han and S. Verdú. Generalizing the Fano inequality. *IEEE Trans. Inform. Theory*, 40(4):1247–1251, 1994.
- [11] T. Kühn. A lower estimate for entropy numbers. *J. Approx. Theory*, 110(1):120–124, 2001.
- [12] E. Lehman and G. Casella. *Theory of point estimation*. Springer-Verlag, New York, NY, 1998.
- [13] K. Lounici. Generalized mirror averaging and d -convex aggregation. *Math. Methods Stat.*, 16(3):246–259, 2007.
- [14] G. Raskutti, M. Wainwright, and B. Yu. Lower bounds on minimax rates for nonparametric regression with additive sparsity and smoothness. In *Proc. Adv. in Neural Proc. Systems (NIPS)*, Vancouver, BC, Dec. 2009.
- [15] G. Raskutti, M.J. Wainwright, and B. Yu. Minimax rates of estimation for high-dimensional linear regression over ℓ_q -balls. *IEEE Trans. Inform. Theory*, 57(10):6976–6994, 2011.
- [16] P. Rigollet and A. Tsybakov. Exponential screening and optimal rates of sparse estimation. *Ann. Stat.*, 39(2):731–771, 2011.
- [17] S. Sarvotham, D. Baron, and R. Baraniuk. Measurements vs. bits: Compressed sensing meets information theory. In *Proc. Allerton Conf. Communication, Control, and Computing*, Monticello, IL, Sept. 2006.
- [18] J. Tropp. User-friendly tail bounds for sums of random matrices. *Found. Comput. Math.*, 12(4):389–434, 2012.
- [19] F. Ye and C. Zhang. Rate minimaxity of the Lasso and Dantzig selector for the ℓ_q loss in ℓ_r balls. *J. Machine Learning Research*, 11:3519–3540, 2010.